

Computer Vision

Computer Vision problems:

- image classification
- object detection
- Neural style Transfer

Deep learning on large images.

small image: $64 \times 64 \times 3 = 12288$ feasible input for Deep neural network.

Actual images in daily life: $1000 \times 1000 \times 3 = 3 \text{ million}$
(this is a relatively small image in actual life).

input size: $(3M, m)$

layer 1: $(n_1, 3M)$ This is getting unpractical.
(too many weights in the network).

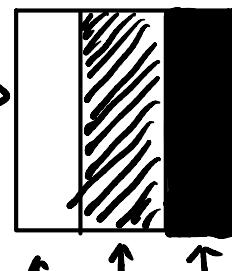
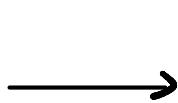
* Solution: convolution.

Edge Detection

vertical edge detection:

filter:
(kernel)

1	0	-1
1	0	-1
1	0	-1



This finds shapes in images where left is bright and right is dark (an edge).

Instead of specifying the filter explicitly, we can let computer learn to detect edges

w ₁	w ₂	w ₃
w ₄	w ₅	w ₆
w ₇	w ₈	w ₉

padding =

Applying a filter on an image has 2 down effects:

- Shrinking output (3x3 filter on 6x6 image, resulting in a 4x4 image)
- * most importantly: throw away information from edge: **edge pixels are used less frequently while applying filter.**

padding adds a layer of blank (zero) pixels outside of the image, solving above 2 problems.

$$\begin{matrix} n \times n \\ \text{image} \end{matrix} \longrightarrow n-f+2p+1 \times n-f+2p+1$$

where f = filter length
 p = padding

valid and same convolutions:

valid: Size of image shrinks after convolution

same: keep output size the same as input size by padding.

strided convolutions:

filters don't have to move across images pixel by pixel. Instead, it can be move by any size for each step =

stride = \geq : move \geq pixels at time.

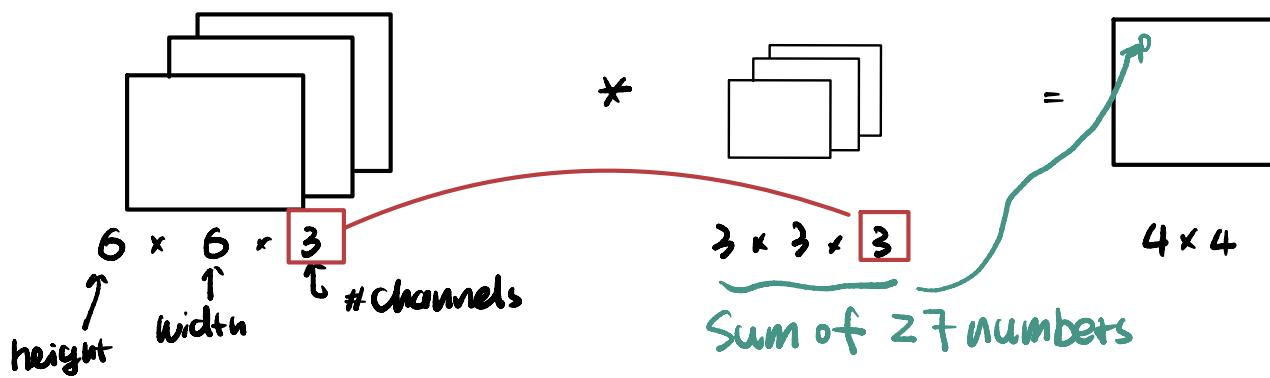
$$n \times n \longrightarrow \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor \times \dots$$

strides, $f \times f$ filter, padding p .

One layer of convolutional network

convolutions over volumes :

common example: convolutions on RGB images



Example: R

1	0	-1
1	0	-1
1	0	-1

G

0	0	0
0	0	0
0	0	0

B

0	0	0
0	0	0
0	0	0

Extract vertical lines in Red channel

R

1	0	-1
1	0	-1
1	0	-1

G

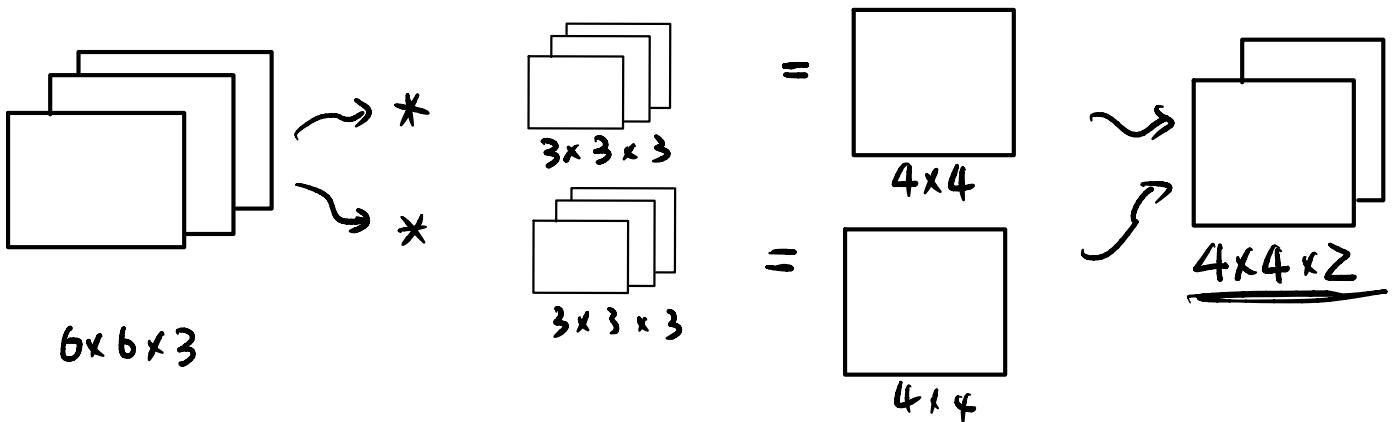
1	0	-1
1	0	-1
1	0	-1

B

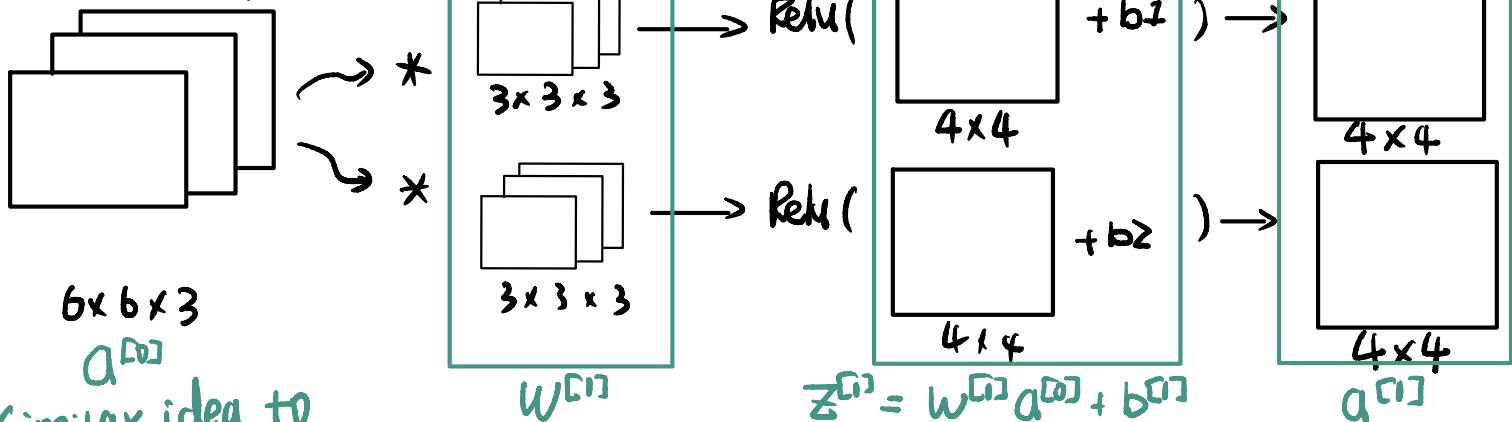
1	0	-1
1	0	-1
1	0	-1

Extract vertical features in all channels

Multiple filters:



Example of a layer



* Similar idea to standard neural network (non-convoluted NN)

Analysis: if you have 10 filters that are 3×3 in one layer of NN,

$$\# \text{parameters} = (\underbrace{3 \times 3 \times 3}_{w} + \underbrace{1}_{b}) \times 10 = \underline{\underline{280}}$$

way less weights than non-convoluted neural network.

Summary of notation:

$f^{[l]}$: filter size

input: $n_H^{[l-1]} \times n_W^{[l-1]} \times n_C^{[l-1]}$

$p^{[l]}$: padding

output: $n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

$s^{[l]}$: stride

$A^{[l]}: \underbrace{m \times n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}}_{\text{1st index}}$

$n_C^{[l]}$: number of filters

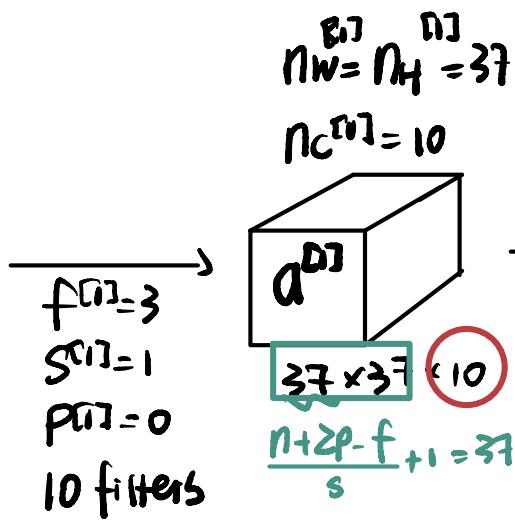
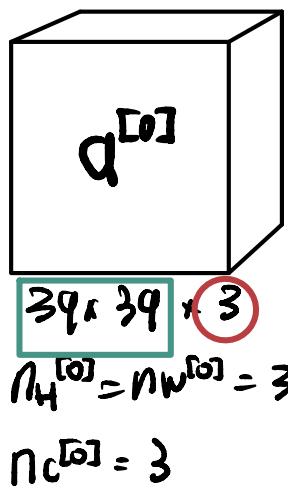
Each filter is: $f^{[l]} \times f^{[l]} \times n_C^{[l-1]}$

Activations: $a^{[l]} \rightarrow n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]} \rightarrow \underbrace{n_H^{[l]}}_{= \frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1} \times n_W^{[l]} \times n_C^{[l]}$

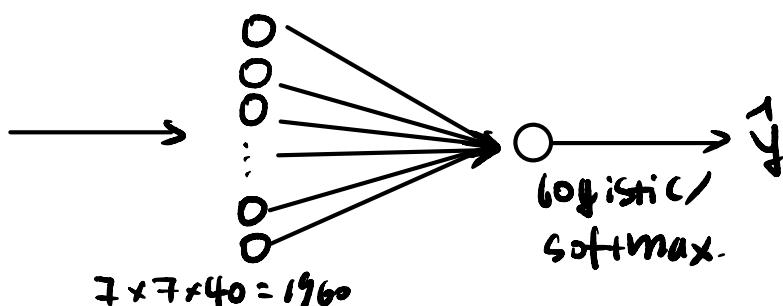
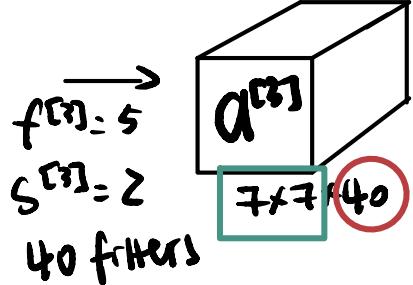
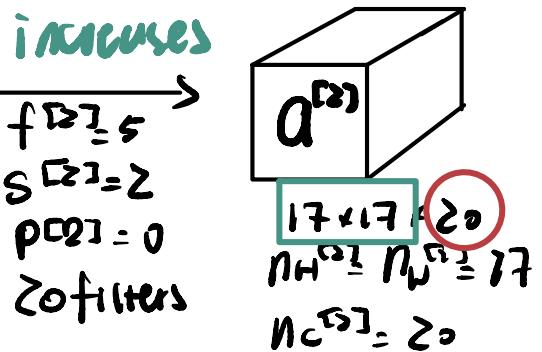
Weights: $f^{[l]} \times f^{[l]} \times n_C^{[l-1]} \times \underbrace{n_C^{[l]}}_{\# \text{filters at this layer}}$

bias: $n_C^{[l]} = (1, 1, 1, \dots, 1, n_C^{[l]})$

Example ConvNet



Notice that convolution size get trim down, while Depth normally increases



* Types of layer in convolutional network.

- Convolution (Conv)
- Pooling (Pool)
- Fully Connected (FC)

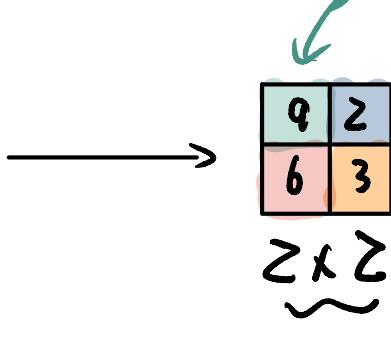
Pooling Layers

Max pooling

1	3	2	1
2	9	1	1
1	3	2	3
5	6	1	2

4×4

e.g. $q = \max(1, 2, 3, 9)$



Same output size as:

$f = 2$
 $s = 2$

hyperparameters

There are other types of pools:

e.g. average pooling (not used as frequently)

Summary:

Hyperparameter:

f : filter size

s : stride

Max or average pooling

Common values:

$f=2$ $s=2$

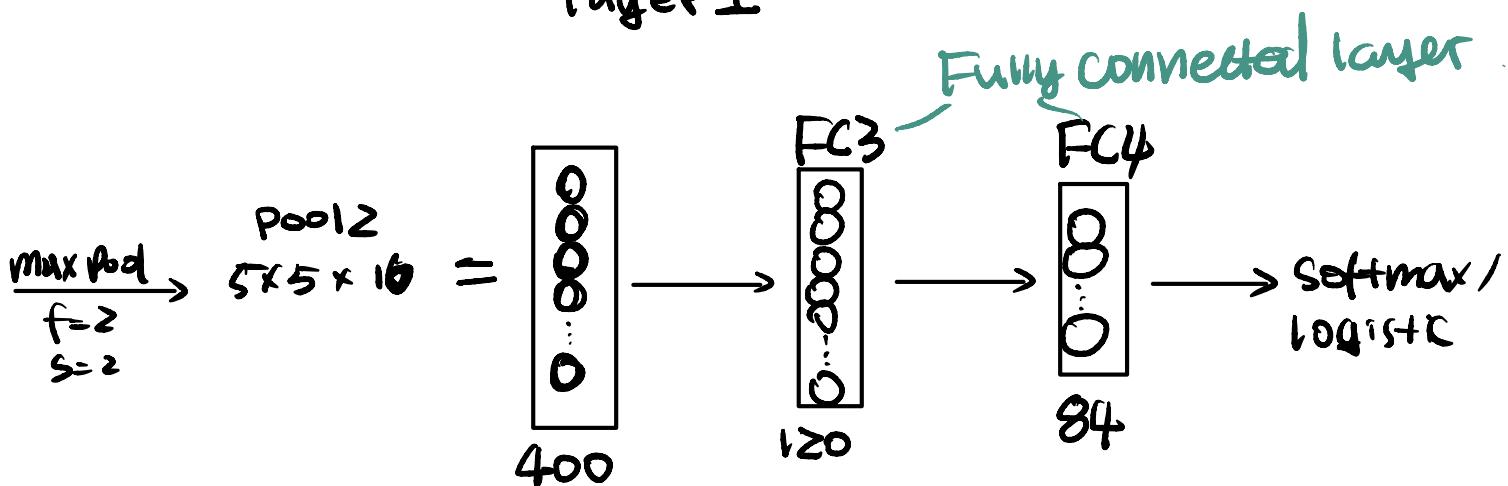
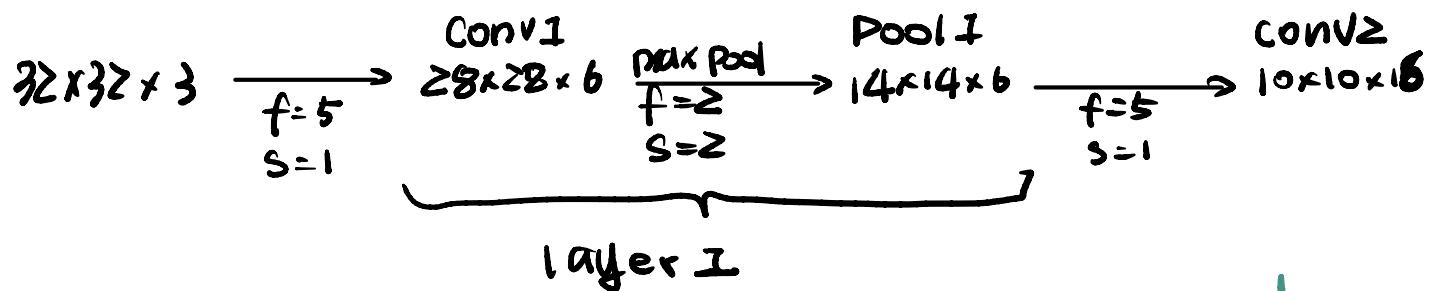
$f=3$ $s=2$

* No parameters to learn.

CNN Example

Example: LeNet-5

p.s. conv1 and pool1 are considered together as layer1, since pool1 has no weights.
no weights
not considered
as independent layer



Common Pattern: H_w, N_w increase

N_C decrease.

Why Convolutions?

- parameter sharing :

e.g. a horizontal line detector can work well across the whole image with only a small number of weights.

- Sparsity of Connections:

in each layer, each output value depends only on a small number of inputs.

Avoid variance given the small amount of weights CNN has.

Training set : $(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})$

Cost : $J = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$

Use gradient descent to optimize parameters to reduce J